

Classification and Representation

Classification

If we use linear regression, also $y = 0$ or 1 .

$h(x)$ can be >1 or <0

So we consider using logistic regression: $0 < h(x) < 1$

To attempt classification, one method is to use linear regression and map all predictions greater than 0.5 as a 1 and all less than 0.5 as a 0. However, this method doesn't work well because classification is not actually a linear function.

The classification problem is just like the regression problem, except that the values we now want to predict take on only a small number of discrete values. For now, we will focus on the **binary classification problem** in which y can take on only two values, 0 and 1. (Most of what we say here will also generalize to the multiple-class case.) For instance, if we are trying to build a spam classifier for email, then $x^{(i)}$ may be some features of a piece of email, and y may be 1 if it is a piece of spam mail, and 0 otherwise. Hence, $y \in \{0,1\}$. 0 is also called the negative class, and 1 the positive class, and they are sometimes also denoted by the symbols "-" and "+" respectively. Given $x^{(i)}$, the corresponding $y^{(i)}$ is also called the label for the training example.

Hypothesis Representation

Logistic Regression Model

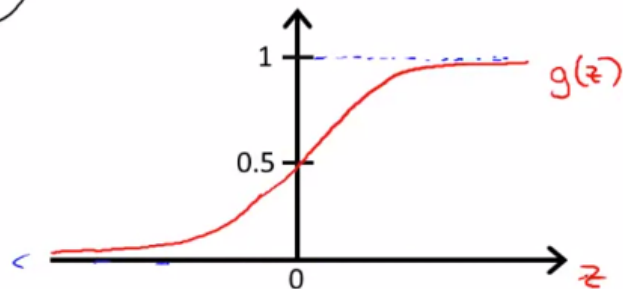
Want $0 \leq h_{\theta}(x) \leq 1$

$$h_{\theta}(x) = g(\theta^T x)$$

$$\rightarrow g(z) = \frac{1}{1 + e^{-z}}$$

Sigmoid function
Logistic function

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$



Interpretation of Hypothesis Output

$h_{\theta}(x)$

$h_{\theta}(x)$ = estimated probability that $y = 1$ on input x ←

Example: If $x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \leftarrow \\ \text{tumorSize} \leftarrow \end{bmatrix}$

$$h_{\theta}(x) = 0.7 \quad y = 1$$

Tell patient that 70% chance of tumor being malignant

$$h_{\theta}(x) = P(y=1|x;\theta)$$

$$y = 0 \text{ or } 1$$

“probability that $y = 1$, given x ,
parameterized by θ ”

$$\begin{aligned} \rightarrow P(y=0|x;\theta) + P(y=1|x;\theta) &= 1 \\ P(\overline{y=0}|x;\theta) &= 1 - P(y=1|x;\theta) \end{aligned}$$

We could approach the classification problem ignoring the fact that y is discrete-valued, and use our old linear regression algorithm to try to predict y given x . However, it is easy to construct examples where this method performs very poorly. Intuitively, it also doesn't make sense for $h_{\theta}(x)$ to take values larger than 1 or smaller than 0 when we know that $y \in \{0, 1\}$. To fix this, let's change the form for our hypotheses $h_{\theta}(x)$ to satisfy $0 \leq h_{\theta}(x) \leq 1$. This is accomplished by plugging $\theta^T x$ into the Logistic Function.

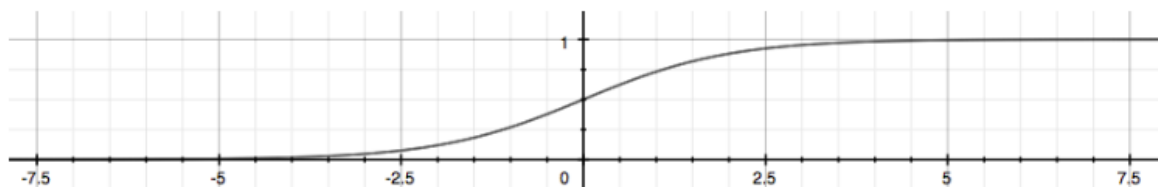
Our new form uses the "Sigmoid Function," also called the "Logistic Function":

$$h_{\theta}(x) = g(\theta^T x)$$

$$z = \theta^T x$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

The following image shows us what the sigmoid function looks like:



The function $g(z)$, shown here, maps any real number to the $(0, 1)$ interval, making it useful for transforming an arbitrary-valued function into a function better suited for classification.

$h_{\theta}(x)$ will give us the **probability** that our output is 1. For example, $h_{\theta}(x) = 0.7$ gives us a probability of 70% that our output is 1. Our probability that our prediction is 0 is just the complement of our probability that it is 1 (e.g. if probability that it is 1 is 70%, then the probability that it is 0 is 30%).

$$\begin{aligned} h_{\theta}(x) &= P(y=1|x;\theta) = 1 - P(y=0|x;\theta) \\ P(y=0|x;\theta) + P(y=1|x;\theta) &= 1 \end{aligned}$$

Decision Boundaries

Logistic regression

$$\rightarrow h_{\theta}(x) = g(\theta^T x) = \underline{p(y=1|x;\theta)}$$

$$\rightarrow g(z) = \frac{1}{1+e^{-z}}$$

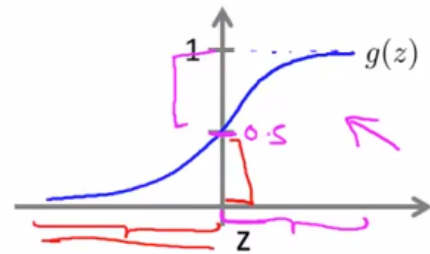
Suppose predict "y = 1" if $h_{\theta}(x) \geq 0.5$

$$\rightarrow \theta^T x \geq 0$$

predict "y = 0" if $h_{\theta}(x) < 0.5$

$$h_{\theta}(x) = g(\theta^T x)$$

$$\rightarrow \theta^T x < 0$$



$$g(z) \geq 0.5$$

when $z \geq 0$

$$h_{\theta}(x) = g(\theta^T x) \geq 0.5$$

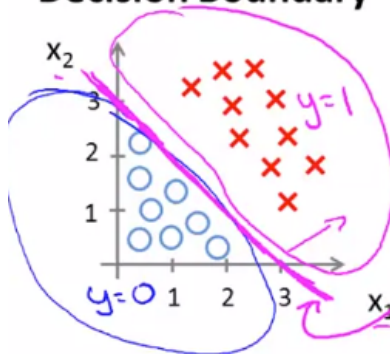
whenever $\theta^T x \geq 0$

$$\uparrow$$

$$z$$

$$g(z) < 0.5$$

Decision Boundary



$$\rightarrow h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

Decision boundary

Predict "y = 1" if $-3 + x_1 + x_2 \geq 0$

$$\theta^T x$$

$$\rightarrow x_1 + x_2 \geq 3$$

$$x_1, x_2 \rightarrow h_{\theta}(x) = 0.5$$

$$x_1 + x_2 = 3$$

$$x_1 + x_2 < 3 \rightarrow y = 0$$

$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix} \leftarrow$$

The decision boundary is a property of the hypothesis. Even if we take away the data set, this decision boundary and the region where we predict $y=1$ versus $y=0$ that's a property of the hypothesis and of the parameters of the hypothesis and not a property of the data set.

In order to get our discrete 0 or 1 classification, we can translate the output of the hypothesis function as follows:

$$\begin{aligned}h_{\theta}(x) &\geq 0.5 \rightarrow y = 1 \\h_{\theta}(x) &< 0.5 \rightarrow y = 0\end{aligned}$$

The way our logistic function g behaves is that when its input is greater than or equal to zero, its output is greater than or equal to 0.5:

$$\begin{aligned}g(z) &\geq 0.5 \\ \text{when } z &\geq 0\end{aligned}$$

Remember.

$$\begin{aligned}z = 0, e^0 = 1 &\Rightarrow g(z) = 1/2 \\ z \rightarrow \infty, e^{-\infty} \rightarrow 0 &\Rightarrow g(z) = 1 \\ z \rightarrow -\infty, e^{\infty} \rightarrow \infty &\Rightarrow g(z) = 0\end{aligned}$$

So if our input to g is $\theta^T X$, then that means:

$$\begin{aligned}h_{\theta}(x) &= g(\theta^T x) \geq 0.5 \\ \text{when } \theta^T x &\geq 0\end{aligned}$$

From these statements we can now say:

$$\begin{aligned}\theta^T x &\geq 0 \Rightarrow y = 1 \\ \theta^T x &< 0 \Rightarrow y = 0\end{aligned}$$

The **decision boundary** is the line that separates the area where $y = 0$ and where $y = 1$. It is created by our hypothesis function.

Example:

$$\begin{aligned}\theta &= \begin{bmatrix} 5 \\ -1 \\ 0 \end{bmatrix} \\ y &= 1 \text{ if } 5 + (-1)x_1 + 0x_2 \geq 0 \\ 5 - x_1 &\geq 0 \\ -x_1 &\geq -5 \\ x_1 &\leq 5\end{aligned}$$

In this case, our decision boundary is a straight vertical line placed on the graph where $x_1 = 5$, and everything to the left of that denotes $y = 1$, while everything to the right denotes $y = 0$.

Again, the input to the sigmoid function $g(z)$ (e.g. $\theta^T X$) doesn't need to be linear, and could be a function that describes a circle (e.g. $z = \theta_0 + \theta_1 x_1^2 + \theta_2 x_2^2$) or any shape to fit our data.